# Predictive Modeling: Principles and Practices

Rick Hefner, Dean Caccavo
Northrop Grumman Corporation

Philip Paul, Rasheed Baqai
Unlimited Innovation, Inc.

**NDIA Systems Engineering Conference**
20-23 October 2008

# Background

- **Predictive modeling relies on historical program performance data (predictive analytics) in conjunction with a forecasting algorithm model to predict future outcomes**

  - Ranges from simple extrapolation techniques to sophisticated Neural Network based models

- **This presentation will discuss the principles of predictive modeling, outline the fundamental methods and tools, and present typical results from applying these techniques to project performance**

# Agenda

- ✓ **What is Predictive Analysis?**
- ▪ **Recent Trends**
- ▪ **Application to Program Performance**
- ▪ **Pilot Results and Feedback**
- ▪ **Summary**

# What is Predictive Analysis?

- ***Could this network packet be from a virus attack?***
  - Predict likelihood of the network packet pattern
  - ➔ **Anomaly detection (outlier detection)**
  - Similar questions:
    - Are the hospital lab results normal (Adverse drug effect detection)
    - Is this credit transaction fraudulent? (fraud detection)

- ***Will this student go to college?***
  - Based on Gender, ParentIncome, ParentEncouragement, IQ, etc.
  - E.g., if ParentEncouragement=Yes and IQ>100, College=Yes
  - ➔ **Classification (prediction)**
  - Similar questions:
    - Is this a spam email? (spam filtering)
    - Recognition of hand-written letters (pen recognition)

- ***What is the person's age?***
  - Based on Hobby, MaritalStatus, NumberOfChildren, Income, HouseOwnership, NumberOfCars, …
  - E.g., If MaritalStatus=Yes, Age = 20+4*NumberOfChildren+0.0001*Income+…
  - ➔ **Regression (prediction)**

## Agenda

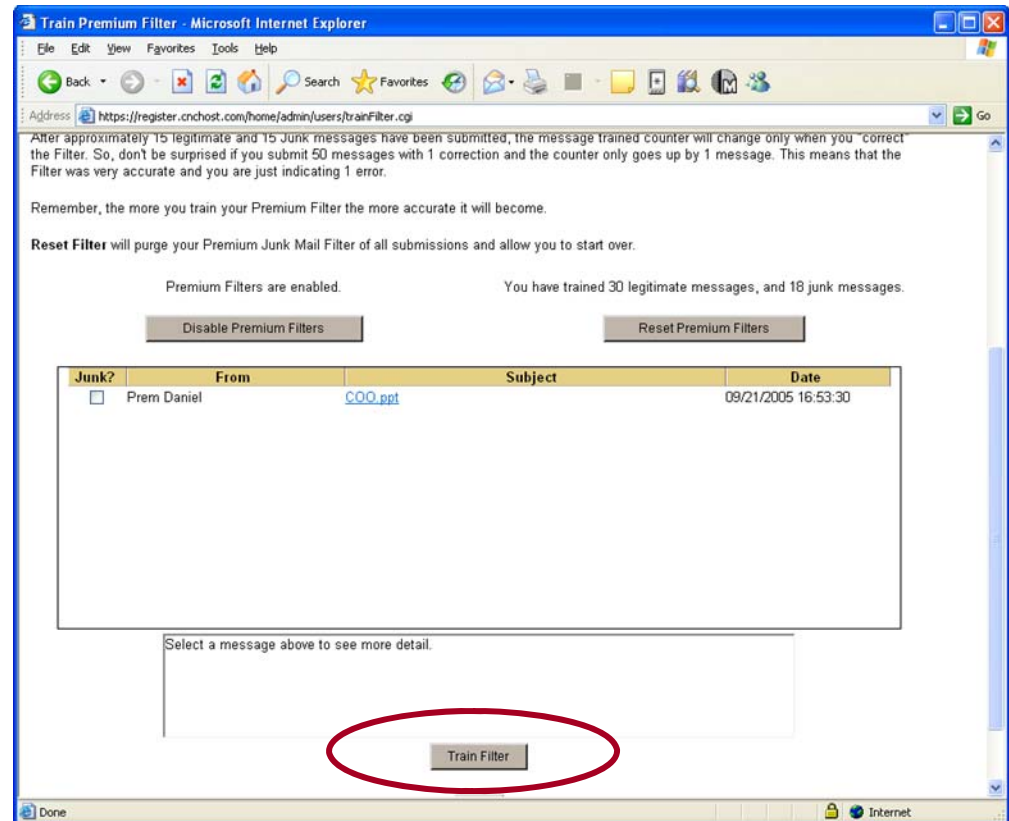- **What is Predictive Analysis?**
- ✓ **Recent Trends**
- **Application to Program Performance**
- **Pilot Results and Feedback**
- **Summary**

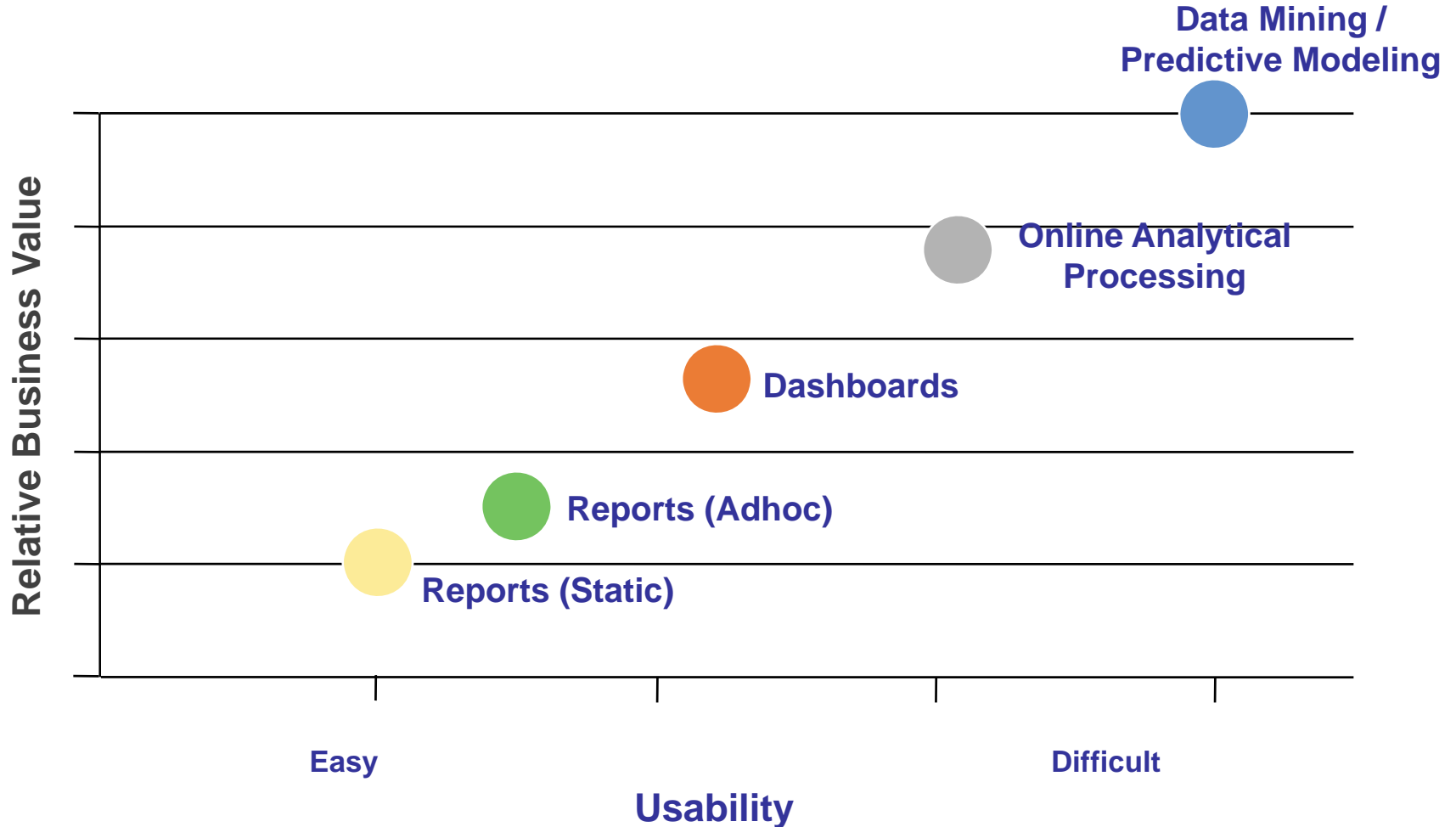# **Predictive Analysis Trends** – Adoption is on the rise

- **Predictive Analysis is becoming more prevalent and integrated in business applications**
  - o *Example: Disease management and evidence based care, based on historical diagnosis and procedure codes of patients*
  - o *Example: E-Mail filtering using predictive analysis*

- **Predictive Analysis algorithms are being integrated into existing databases, data mining tools**
  - o *Example: Microsoft SQL Server 2005 has predictive analysis algorithms*

Example:
Premium predictive analysis based filtering on e-mail, available to any e-mail user

# **Predictive Analysis Trends** – Tools are becoming easier to use



**Data Mining / Predictive Modeling**

**Online Analytical Processing**

**Dashboards**

**Reports (Adhoc)**

**Reports (Static)**

Relative Business Value

Easy

Difficult

**Usability**

# **Predictive Analysis Trends** – Model development is more structured

**Define a Model**

**Train the Model**

Training Data

**Test the Model**
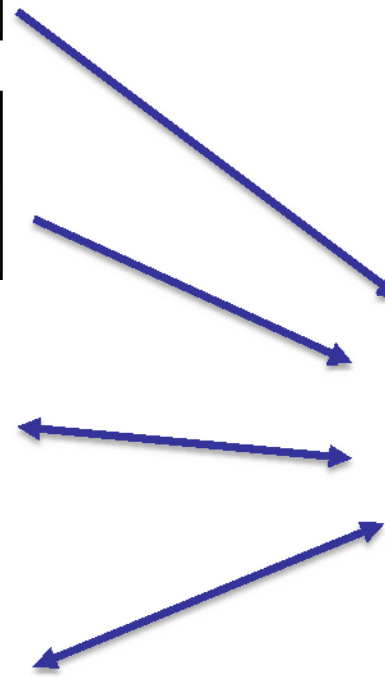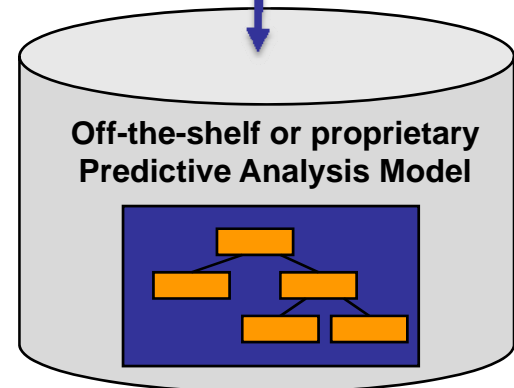
Test Data

**Prediction using the Model**

Prediction Input Data

**Off-the-shelf or Proprietary Predictive Analysis Engine**

Third Party Predictive Analysis tools

**Off-the-shelf or proprietary Predictive Analysis Model**

- Executive understanding of the creation, training and testing of the model is critical to success
- The Model gets more powerful and accurate as the volume of data fed into the model increases

## **Predictive Analysis Trends** – Algorithms are available for use

| Decision Trees | Naïve Bayesian | Clustering | Sequential Clustering | Time Series | Association rules | Neural Network | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---|
| 1 | 2 | 2 | 2 | | 2 | 1 | **Classification** |
| 1 | 2 | 2 | 2 | | | 1 | **Regression** |
| | | 1 | 1 | | | 2 | **Segmentation** |
| 1 | 1 | 2 | 2 | | 1 | 1 | **Association Analysis** |
| | | 1 | 1 | | | 2 | **Anomaly Detect.** |
| | | | 1 | | | | **Sequential  Analysis** |
| | | | | 1 | | | **Time series** |

**1 - First Choice**       2 - Second Choice

## Data Mining Vendors & Tools

- SAS (Enterprise Miner)
- IBM (DB2 Intelligent Miner)
- Oracle (ODM option to Oracle 10g)
- SPSS (Clementine)
- Insightful (Insightful Miner)
- KXEN (Analytic Framework)
- Prudsys (Discoverer and its family)
- Microsoft (SQL Server 2005)
- Angoss (KnowledgeServer and its family)
- DBMiner (DBMiner)
- Many others

# Agenda

- ■ **What is Predictive Analysis?**
- ■ **Recent Trends**
- ✓ **Application to Program Performance**
- ■ **Pilot Results and Feedback**
- ■ **Summary**

# Mission Assurance Continuum

| Program Performance Oversight | Program Analysis Reporting | Predictive Program Health |
|---|---|---|
| **Industry Minimum** | **Industry Best Practice** | **Industry Innovators** |
| Proactive Program Management Program Portfolio Management | Reports based on current and passed performance data of portfolio programs, programs, and subcontract reports | Predictive Analysis based on Program Performance Modeling |

### Approach and Scope

| | | • Self reported program metrics, organizational data, personnel data and customer reported metrics collected at regular intervals |
|---|---|---|
| • Self reported Program Portfolio includes critical and high visibility programs | • Self Reported Program metrics collected periodically and at specific program milestones | • Predictive models developed using historical data (leading indicators rationalized) |
| • Standard Program Management Metrics collected on a periodic basis | • Reporting analysis performed as needed | • Models validated against historical data |

### Infrastructure and Breadth

| | | • Holistic enterprise wide approach to program execution |
|---|---|---|
| • Program data maintained by individual programs | • Program data collected periodically into an enterprise-wide program management repository | • Models continually refined using current program performance data |
| • Summary information provided to enterprise repository | • Program, Enterprise and Subcontracts performance reports available | • Sophisticated predictive measures provided to programs and enterprise |

### Data Requirements

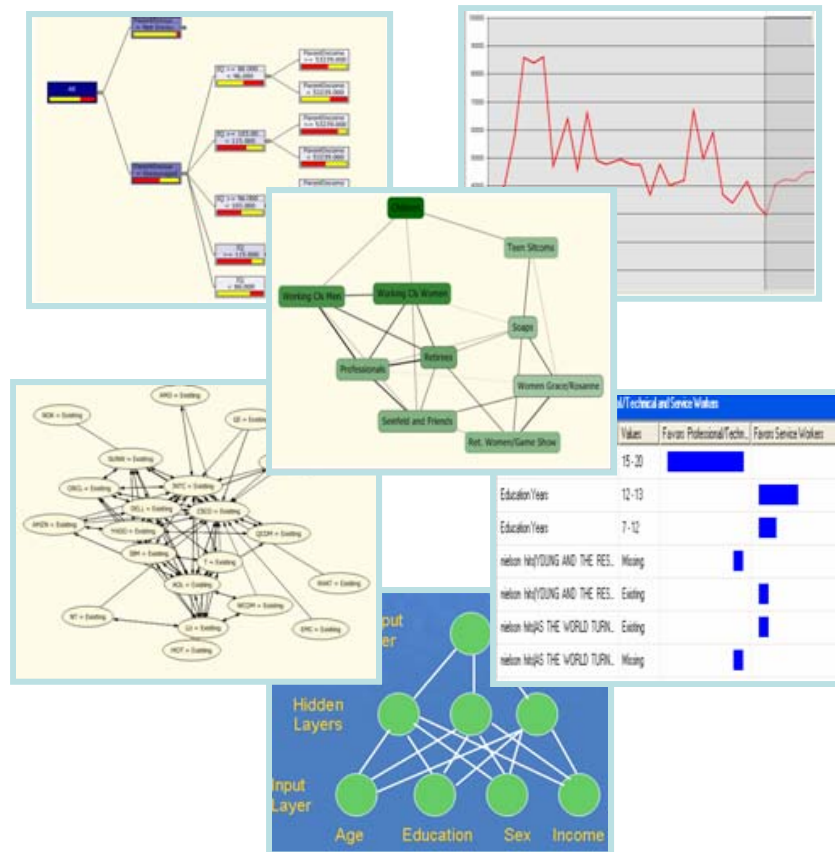| • Very few metrics collected from programs | • 25 – 100 metrics collected from programs | • 50 – 75 metrics collected from programs and refined to include only the few relevant metrics |
|---|---|---|
| • Key program metrics (cost performance, schedule performance, technical performance, CPI, SPI etc.) | • Key program metrics collected at all specified Program Milestones. | • Adaptive approach to qualitative and quantitative performance indicators |
| • Standardized program taxonomy information like customer, contract type | | • Direct and Indirect metrics collected for the programs; qualitative information is mined |
| | | • Proactive responses based on predictive analysis of ongoing and historical performance |

## Overarching Objectives for Predictive Modeling

- **Provide program management staff with Predictive Models to "test-their-gut" against enterprise experience data before making strategic program decisions**

- **Develop Predictive Models that provide insight into identifying "headlight metrics" that influence Schedule and Cost realism during program execution**

- **Leverage existing enterprise information to develop Predictive Models for programs**

- **Ensure that models are extensible and automatically calibrated with additional data from the program and enterprise**

# Potential Areas for Predictive Analysis

## Potential Predictive Analysis Models for Program Management and Subcontractor Management

- Schedule Risk at WBS level based on past performance
- Cost Risk at WBS level based on past performance
- Technical Risk at WBS level based on past performance
- Spending and staffing profile for the program life cycle
- Subcontractor risk profile based on past performance
- Sub-tier quality at subcontract and WBS level
- Defect/Aberrations for the program life cycle
- Mission Assurance models based on program category



## Predictive Analysis Algorithms

- Decision Trees
- Naïve Bayesian
- Clustering
- Sequence Clustering
- Association Rules
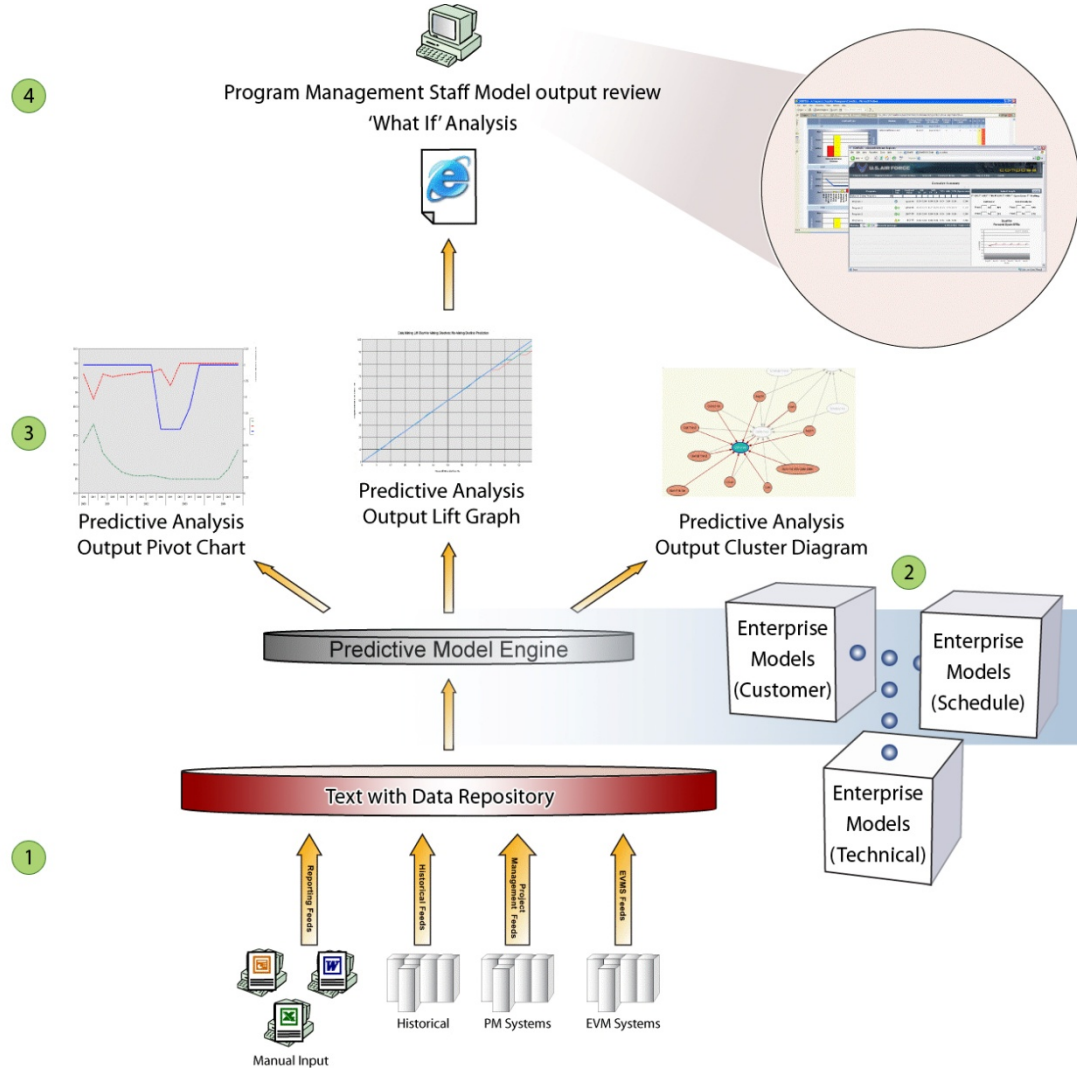- Neural Network
- Time Series
- Custom Model

# Predictive Analysis High Level CONOPS

1) Enterprise data is mined and analyzed

2) Enterprise models are defined by Analysts

3) Enterprise model outputs are defined by Analysts and customized by PM staff

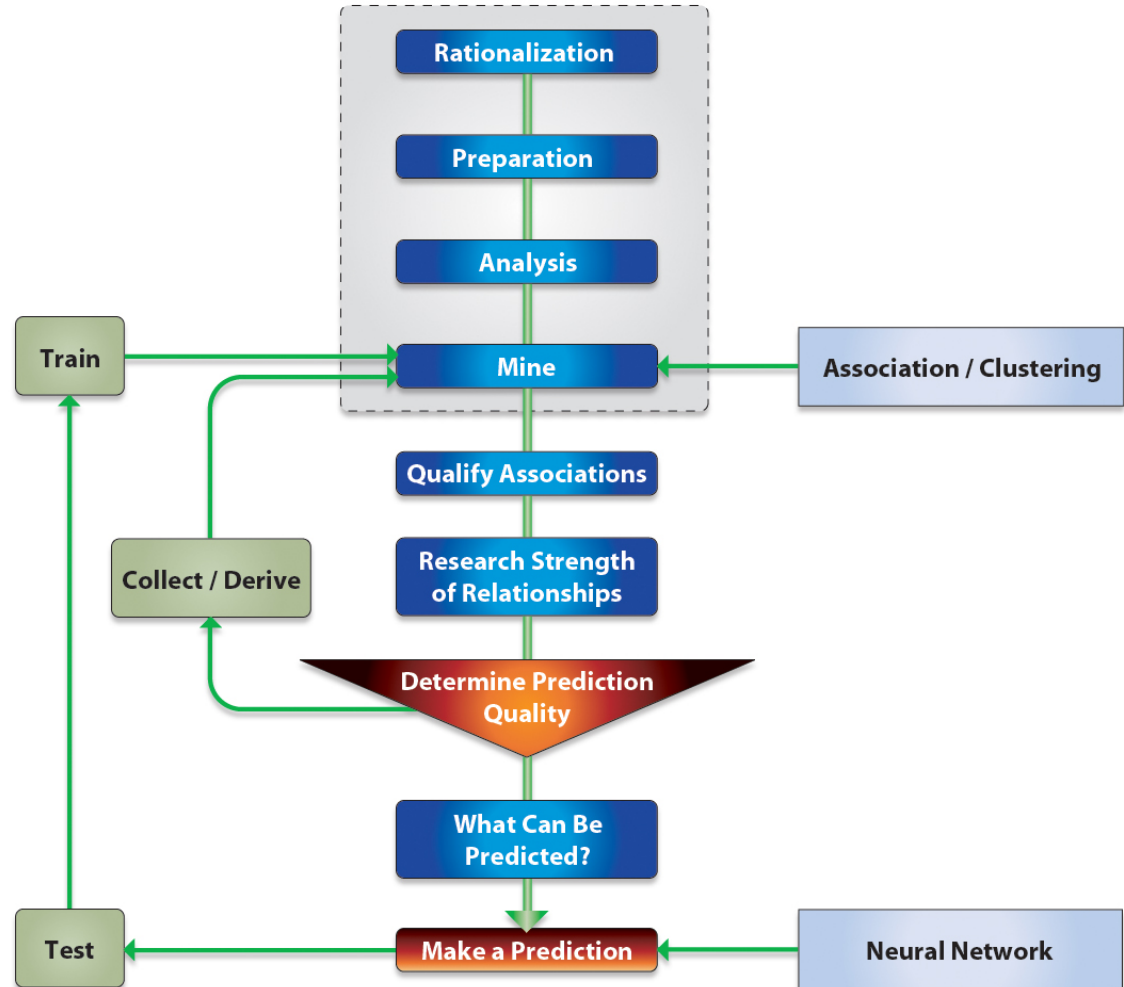4) PM staff use models interactively

**Key Benefit:**
Leverages enterprise experience data and sophisticated algorithms into predictive models for cost and schedule realism checks during program execution

Program Management Staff Model output review
'What If' Analysis

Predictive Analysis Output Pivot Chart

Predictive Analysis Output Lift Graph

Predictive Analysis Output Cluster Diagram

Enterprise Models (Customer)

Enterprise Models (Schedule)

Enterprise Models (Technical)

Predictive Model Engine

Text with Data Repository

Reporting Feeds

Historical Feeds

Project Management Feeds

EVMS Feeds

Manual Input

Historical

PM Systems

EVM Systems

15

# The Predictive Modeling Process

- **Explore the Data**
- **Understand Data Relationships**
- **Derive/Enhance the Data**
- **Use the Data to Predict**
- **Train the Model**

# What can be Predicted with Reasonable Accuracy?

|  | **Limited Number of Programs** | **Enterprise Experience** |
|---|---|---|
| **Large volume of historical data** | ■ Likelihood or return to acceptable performance<br>■ Predictive Program Performance<br><br>**1** | ■ Quadrant 2 predictions<br>■ Quadrant 3 predictions<br>■ Early warning "headlight indicators"<br>■ Higher accuracy based on enterprise experience  **3** |
| **Limited Historical data** | | ■ Cost, schedule realism<br>■ Phase realism<br>■ WBS Accuracy<br><br>**2** |

**Program Lifecycle Stage**

**Low**                      **High**

**Volume of "Like" Programs**

# Agenda

- ■ **Background**
- ■ **Industry Trends**
- ■ **Application to Program Performance**
- ✓ **Pilot Results and Feedback**
- ■ **Summary**

# Predictive Modeling Pilot Objectives

- **Provide program management staff with Predictive Models to "test-their-gut" against enterprise experience data before making strategic program decisions**

- **Develop Predictive Models that provide insight into identifying "headlight metrics" that influence Schedule and Cost realism during program execution**

- **Leverage existing enterprise information to develop Predictive Models for programs**

- **Ensure that models are extensible and automatically calibrated with additional data from the program and enterprise**

# Pilot Approach

- Analyze and rationalize the available enterprise data
    - Enterprise Level Office of Cost Estimation and Risk Assessment (OCERA) data
    - Division Level Stoplight Program data
    - Program Level Program Review Authority (PRA) data for relevant programs
- Develop predictive modeling approach to provide schedule and cost measures during program execution phase
- Develop preliminary predictive models using appropriate algorithms and mining existing enterprise data
    - Mining – Clustering, Decision Trees and Naïve Bayesian Algorithms
    - Predictions – Neural Network, Bayesian Algorithms and Clustering
- Get Pilot participation from three representative program types:
    - Large Scale System Integration Low Rate Initial Production program
    - Medium Sized Software program
    - Small IT System (Software and Hardware) program

**Key Benefit: Leverages enterprise experience data and sophisticated algorithms into predictive models for use during program execution**

# Data analyzed for developing preliminary models

| Data | Stoplight | OCERA | PRA |
|---|---|---|---|
| Data Period | 2.5 years | 5 – 6 years | Past 4 months |
| Frequency | Quarterly/Some older data is monthly | Major milestones or annually | Monthly |
| Breadth and depth of data | Monthly snapshot of key metrics | Very deep, very broad, with significant contextual information | Very deep, mostly snapshot without significant contextual information |
| Approximate number of data elements | ~ 20 | ~ 70 key attributes | ~40 key attributes |

**Analyzed enterprise level (OCERA), division level (Stoplight) and program level (PRA) data**

# Some Actual Data Types Used to Develop Predictive Model Relationships

## Program Data

- Contract Type
  - CPAF, FFP, CPFF
- Type of Program
- Period of Performance
- Number of Milestones
- Number of sub-contractors
  - Subcontract value
  - Subcontract performance
- Total Value
- Annual Sales
- Number of incremental deliveries
- Average staff count
- SPI, CPI
- EAC, BAC
- Number of EAC changes
- Number of ECR/ECP
- Defects
  - Injection by phase
  - Occurrence by phase
- Skills Data
- Program Review Data
- Project Initiation Review Data

## Program Self Assessment

- Monthly Ratings
  - Schedule
  - Technical
  - Cost
  - Mission Assurance
  - Management
  - Process

## External Data

  - CPARS
  - Customer satisfaction data
  - Award Fees

## Milestone Data

- Milestones
  - Proposal
  - Contract Startup
  - SRR
  - SDR
  - Software Specification Review
  - PDR
  - CDR
  - Test Readiness Review
  - Completion

## Other Data

- Action Item Data
- Organization benchmark data
- SLOC, ESLOC
- Productivity
- Language, Component type, complexity,
- Reuse ratios
- Platform, environment

**Contains Enterprise, Division and Program Data**

# Data Mining Results

| Column Name | Score | Input |
|---|---|---|
| PRA | 0.295 | x |
| Tech | 0.275 | x |
| CustSatis | 0.266 | x |
| ContractType | 0.203 | x |
| Customer | 0.154 | x |
| Cost | 0.121 | x |
| ContractTypeID | 0.113 | x |
| MA | 0.101 | x |
| MonthAndYear | 0.093 | x |
| TypeOfWork | 0.084 | x |
| EVM | 0.078 | x |
| Organization | 0.075 | x |
| EVMReq | 0.073 | x |
| Supl | 0.068 | x |
| WorkTypeID | 0.062 | x |
| UnconPrecon | 0.058 | x |
| T2N | 0.034 | |
| Proc1 | 0.030 | |
| CashFlowDSR | 0.023 | |
| POPBegins | 0.000 | |

**Prediction Measures**
- Schedule
- Cost



- **The mining showed that out of the over 125 metrics and measures some are leading indicators and are more important than others in influencing cost and schedule**
- **While it cannot be proved to be conclusive with the limited data that was used, the trends were definite**

# Derivation of Data & Data Relationships

- **Examples of Derived Data**
  - Number of Outstanding Program Issues (with and without recovery dates)
  - Variance in program Cost/Schedule/Technical health from month-to-month
  - Program Cost/Schedule/Technical health trend from month-to-month
  - Variance in VAC from month-to-month taken as a percentage of the current EAC

- **Examples of Discovered Relationships**
  - Schedule Health is a good indicator of program Overall Health recovery
  - Cost and Technical Health are good indicators of program Overall Health decline

**Better understanding of the data allows for organization and enhancement of the dataset**

# Model Development & Calibration

| Model | Calibrated Model |
|---|---|



- Modeling without applied domain knowledge or calibration resulted in lower accuracy
- Association models able to determine relevant data attributes

- Incorporating domain knowledge and calibration into data mining resulted in higher accuracy
- Data relationships are more clearly defined

**Domain knowledge & calibration applied to data mining can enhance the predictive model**
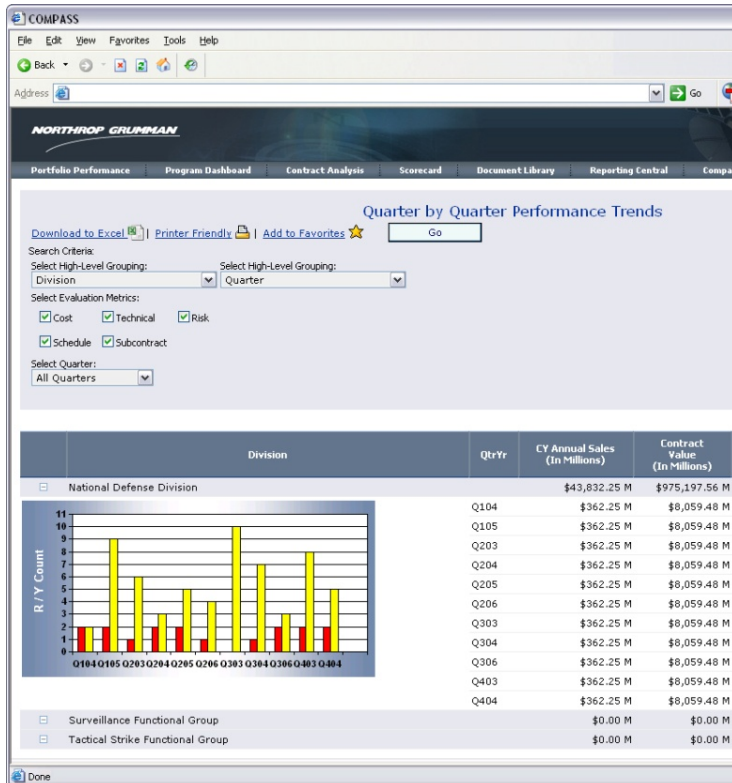
# Typical Results from the Models

Ability for Programs to review the predictive output from multiple models to "test-the-gut" before making strategic program decisions



FICTIONAL DATA

# Typical Results from the Models

Ability for staff to review status and trends across the portfolio of programs, across a variety of categories



**FICTIONAL DATA**

## Agenda

- **What is Predictive Analysis?**
- **Recent Trends**
- **Application to Program Performance**
- ✓ **Summary**

## **Summary** – Critical success factors

- Executive and Enterprise support and understanding of long-term strategic benefits

- Understanding of the types of data and the correlation between the data

- Understanding of the various constituents in the value chain and the tools/processes for each constituent

- Prototypes or mockups that depict the results of the model

- Sound and robust technical architecture

- Delivery mechanism that shields the complexity of the model from the end users

## More Information

- OLE DB for DM specification
  - http://www.microsoft.com/downloads/detail s.aspx?FamilyID=01005f92-dba1-4fa4-8ba0-af6a19d30217&DisplayLang=en
- Plug-in
  - http://www.msnusers.com/AnalysisServic esDataMining/Documents/Files%2FSQL%2 0Server%20Data%20Mining%20Plug%2DI n%20Algorithms%20%28Beta%202%20%2B%2B%29.zip
  - A white paper, tutorial, and complete sample code for Pair-wise Linear Regression
- SQL Server 2005:
  - www.microsoft.com/sql/2005
- Community:
  - Microsoft.public.sqlserver.datamining
  - Microsoft.private.sqlserver2005.analysisser vices.datamining
  - Groups.msn.com/AnalysisServicesDataMin ing
- msdn.microsoft.com (search "data mining")

- Decision trees (classification/regression):
  - ftp://ftp.research.microsoft.com/users/surajitc /icde99.pdf
  - http://www.research.microsoft.com/research/ pubs/view.aspx?tr_id=81
  - http://research.microsoft.com/~dmax/publicat ions/dmart-final.pdf
- Association rules:
  - Apriori algorithm (see Data Mining concepts and techniques)
- Clustering
  - EM:http://www.research.microsoft.com/script s/pubs/view.asp?TR_ID=MSR-TR-98-35
  - K-means (see Data Mining concepts and techniques)
- Sequence clustering
  - ftp://ftp.research.microsoft.com/pub/tr/tr-2000-18.pdf
- Time series:
  - http://research.microsoft.com/~dmax/publicat ions/dmart-final.pdf
- Neural network
  - Conjugate gradient method (see Data Mining concepts and techniques)
- Naïve Bayesian
  - See Data Mining concepts and techniques

## Contact Information



Rick Hefner, Ph.D.

Northrop Grumman Corporation

(310) 812-7290

rick.hefner@ngc.com